

## 10-701 Introduction to Machine Learning (PhD) Lecture 14: Learning Theory

Leila Wehbe  
Carnegie Mellon University  
Machine Learning Department

Slides based on Tom Mitchell's 10701 Fall 2016 material  
Readings: [TM] chapter 7  
Nina Balcan's notes on generalization guarantees: [http://  
www.cs.cmu.edu/~ninamf/courses/601sp15/sc-2015.pdf](http://www.cs.cmu.edu/~ninamf/courses/601sp15/sc-2015.pdf)

## Announcements

- Project report due Friday
  - Follow the requirements!
- Monday 3/18 and Wednesday 3/20:
  - Class starts at 11!
  - We will do a midterm review

## Computational Learning Theory

- What general laws constrain inductive learning?
- Want theory to relate
  - Number of training examples
  - Complexity of hypothesis space
  - Accuracy to which target function is approximated
  - Manner in which training examples are presented
  - Probability of successful learning

\* See annual Conference on Computational Learning Theory

## Sample Complexity

How many training examples suffice to learn target concept

1. If learner proposes instances as queries to teacher?
  - learner proposes  $x$ , teacher provides  $f(x)$
2. If teacher (who knows  $f(x)$ ) generates training examples?
  - teacher proposes sequence  $\{(x^1, f(x^1)), \dots (x^n, f(x^n))\}$
3. If some random process (e.g., nature) generates instances, and teacher labels them?
  - instances drawn according to  $P(X)$

## Sample Complexity 3

Problem setting:

- Set of instances  $X$
- Set of hypotheses  $H = \{h : X \rightarrow \{0,1\}\}$
- Set of possible target functions  $C = \{c : X \rightarrow \{0,1\}\}$
- Sequence of training instances drawn at random from  $P(X)$   
teacher provides noise-free label  $c(x)$

Learner outputs a hypothesis  $h \in H$  such that

$$h = \arg \min_{h \in H} error_{train}(h)$$

## Example: Learning decision trees

Take  $X = (X_1, \dots, X_n)$  s.t.  $X_i \in \{0,1\}$

$Y_i \in \{0,1\}$

Let  $H$  be the set of decision trees:

$$H = \{h : X \rightarrow Y\}$$

How many possible values of  $X$ ?

How many possible trees?

How many training examples needed to find the right tree?

## Example: Learning decision trees

Take  $X = (X_1, \dots, X_n)$  s.t.  $X_i \in \{0,1\}$

$Y_i \in \{0,1\}$

Let  $H$  be the set of decision trees:

$$H = \{h : X \rightarrow Y\}$$

How many possible values of  $X$ ?  $2^n$

How many possible trees?

How many training examples needed to find the right tree?

## Example: Learning decision trees

Take  $X = (X_1, \dots, X_n)$  s.t.  $X_i \in \{0,1\}$

$Y_i \in \{0,1\}$

Let  $H$  be the set of decision trees:

$$H = \{h : X \rightarrow Y\}$$

How many possible values of  $X$ ?  $2^n$

How many possible trees?  $|H| = 2^{2^n}$

How many training examples needed to find the right tree?

### Example: Learning decision trees

Take  $X = (X_1, \dots, X_n)$  s.t.  $X_i \in \{0,1\}$

$Y_i \in \{0,1\}$

Let  $H$  be the set of decision trees:

$H = \{h : X \rightarrow Y\}$

How many possible values of  $X$ ?  $2^n$

How many possible trees?  $|H| = 2^{2^n}$

How many training examples needed to find the right tree?

### Example: Learning decision trees

Take  $X = (X_1, \dots, X_n)$  s.t.  $X_i \in \{0,1\}$

$Y_i \in \{0,1\}$

Let  $H$  be the set of decision trees:

$H = \{h : X \rightarrow Y\}$

How many possible values of  $X$ ?  $2^n$

How many possible trees?  $|H| = 2^{2^n}$

How many training examples needed to find the right tree?  $2^n$   
(no free lunch)

### Example: Learning decision trees

Take  $X = (X_1, \dots, X_n)$  s.t.  $X_i \in \{0,1\}$

$Y_i \in \{0,1\}$

Let  $H$  be the set of decision trees:

$H = \{h : X \rightarrow Y\}$

How many possible values of  $X$ ?  $2^n$

How many possible trees?  $|H| = 2^{2^n}$

How many training examples needed to find the right tree?  $2^n$   
(no free lunch)

Generalizing beyond training is impossible unless we add assumptions

### Example: Learning decision trees

Take  $X = (X_1, \dots, X_n)$  s.t.  $X_i \in \{0,1\}$

$Y_i \in \{0,1\}$

Let  $H$  be the set of decision trees:

$H = \{h : X \rightarrow Y\}$

How many possible values of  $X$ ?  $2^n$

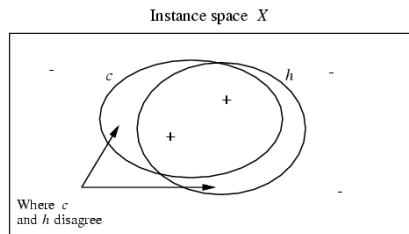
How many possible trees?  $|H| = 2^{2^n}$

How many training examples needed to find the right tree?  $2^n$   
(no free lunch)

Generalizing beyond training is impossible unless we add assumptions

training examples are provided according to distribution  $P(X)$

## True Error of a Hypothesis



The *true error* of  $h$  is the probability that it will misclassify an example drawn at random from  $P(X)$

$$error_{true}(h) \equiv \Pr_{x \sim P(X)} [h(x) \neq c(x)]$$

## Two notions of error

*Training error* of hypothesis  $h$  with respect to target concept  $c$

- How often  $h(x) \neq c(x)$  over training instances  $D$

$$error_{train}(h) \equiv \Pr_{x \in D} [h(x) \neq c(x)] = \frac{1}{|D|} \sum_{x \in D} \delta(h(x) \neq c(x))$$

*True error* of hypothesis  $h$  with respect to  $c$

- How often  $h(x) \neq c(x)$  over future instances drawn at random from  $\mathcal{D}$

$$error_{true}(h) \equiv \Pr_{x \sim P(X)} [h(x) \neq c(x)]$$

training examples  $D$

Probability distribution  $P(X)$

## Overfitting

Consider a hypothesis  $h$  and its

- Error rate over training data:  $error_{train}(h)$
- True error rate over all data:  $error_{true}(h)$

We say  $h$  overfits the training data if

$$error_{true}(h) > error_{train}(h)$$

Amount of overfitting =

$$error_{true}(h) - error_{train}(h)$$

## Overfitting

Consider a hypothesis  $h$  and its

- Error rate over training data:  $error_{train}(h)$
- True error rate over all data:  $error_{true}(h)$

We say  $h$  overfits the training data if

$$error_{true}(h) > error_{train}(h)$$

Amount of overfitting =

$$error_{true}(h) - error_{train}(h)$$

Can we bound  $error_{true}(h)$   
in terms of  $error_{train}(h)$  ??

$$error_{train}(h) \equiv \Pr_{x \in D} [h(x) \neq c(x)] = \frac{1}{|D|} \sum_{x \in D} \delta(h(x) \neq c(x))$$

training  
examples

$$error_{true}(h) \equiv \Pr_{x \sim P(X)} [h(x) \neq c(x)]$$

Probability  
distribution  $P(x)$

if  $D$  was a set of examples drawn from  $P(X)$  and **independent** of  $h$ , then we could use standard statistical confidence intervals to determine that with 95% probability,  $error_{true}(h)$  lies in the interval:

$$error_D(h) \pm 1.96 \sqrt{\frac{error_D(h)(1 - error_D(h))}{n}}$$

but  $D$  is the **training data** for  $h$  ....

## Version Spaces

$$c : X \rightarrow \{0,1\}$$

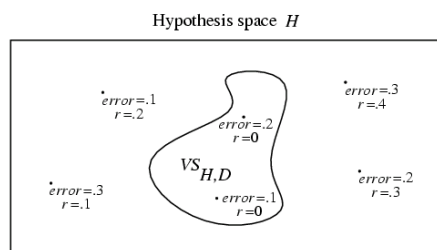
A hypothesis  $h$  is **consistent** with a set of training examples  $D$  of target concept  $c$  if and only if  $h(x) = c(x)$  for each training example  $\langle x, c(x) \rangle$  in  $D$ .

$$Consistent(h, D) \equiv (\forall \langle x, c(x) \rangle \in D) h(x) = c(x)$$

The **version space**,  $VS_{H,D}$ , with respect to hypothesis space  $H$  and training examples  $D$ , is the subset of hypotheses from  $H$  consistent with all training examples in  $D$ .

$$VS_{H,D} \equiv \{h \in H \mid Consistent(h, D)\}$$

## Exhausting the version space



( $r$  = training error,  $error$  = true error)

**Definition:** The version space  $VS_{H,D}$  with respect to training data  $D$  is said to be  **$\epsilon$ -exhausted** if every hypothesis  $h$  in  $VS_{H,D}$  has true error less than  $\epsilon$ .

$$(\forall h \in VS_{H,D}) error_{true}(h) < \epsilon$$

## How many examples will $\epsilon$ -exhaust the version space?

**Theorem:** [Haussler, 1988].

If the hypothesis space  $H$  is finite, and  $D$  is a sequence of  $m \geq 1$  independent random examples of some target concept  $c$ , then for any  $0 \leq \epsilon \leq 1$ , the probability that the version space with respect to  $H$  and  $D$  is not  $\epsilon$ -exhausted (with respect to  $c$ ) is less than

$$|H|e^{-\epsilon m}$$

## How many examples will $\epsilon$ -exhaust the version space?

**Theorem:** [Haussler, 1988].

If the hypothesis space  $H$  is finite, and  $D$  is a sequence of  $m \geq 1$  independent random examples of some target concept  $c$ , then for any  $0 \leq \epsilon \leq 1$ , the probability that the version space with respect to  $H$  and  $D$  is not  $\epsilon$ -exhausted (with respect to  $c$ ) is less than

$$|H|e^{-\epsilon m}$$

Interesting! This bounds the probability that any consistent learner will output a hypothesis  $h$  with  $error(h) \geq \epsilon$

Any(!) learner that outputs a hypothesis consistent with all training examples (i.e., an  $h$  contained in  $VSH,D$ )

## What it means

[Haussler, 1988]: probability that the version space is not  $\epsilon$ -exhausted after  $m$  training examples is at most  $|H|e^{-\epsilon m}$

$$\Pr[(\exists h \in H) s.t. (error_{train}(h) = 0) \wedge (error_{true}(h) > \epsilon)] \leq |H|e^{-\epsilon m}$$

Suppose we want this probability to be at most  $\delta$

1. How many training examples suffice?

$$m \geq \frac{1}{\epsilon} (\ln |H| + \ln(1/\delta))$$

2. If  $error_{train}(h) = 0$  then with probability at least  $(1-\delta)$ :

$$error_{true}(h) \leq \frac{1}{m} (\ln |H| + \ln(1/\delta))$$

## Example: $H$ is Conjunction of up to $N$ Boolean Literals

Consider classification problem  $f: X \rightarrow Y$ :  $m \geq \frac{1}{\epsilon} (\ln |H| + \ln(1/\delta))$

- instances:  $X = (X_1, X_2, X_3, X_4)$  where each  $X_i$  is boolean
- Each hypothesis in  $H$  is a rule of the form:
  - IF  $(X_1, X_2, X_3, X_4) = (0, ?, 1, ?)$ , THEN  $Y=1$ , ELSE  $Y=0$
  - i.e., rules constrain any subset of the  $X_i$

How many training examples  $m$  suffice to assure that with probability at least 0.99, any consistent learner using  $H$  will output a hypothesis with true error at most 0.05?

## Example: $H$ is Conjunction of up to $N$ Boolean Literals

Consider classification problem  $f: X \rightarrow Y$ :  $m \geq \frac{1}{\epsilon} (\ln |H| + \ln(1/\delta))$

- instances:  $X = (X_1, X_2, X_3, X_4)$  where each  $X_i$  is boolean
- Each hypothesis in  $H$  is a rule of the form:
  - IF  $(X_1, X_2, X_3, X_4) = (0, ?, 1, ?)$ , THEN  $Y=1$ , ELSE  $Y=0$
  - i.e., rules constrain any subset of the  $X_i$

How many training examples  $m$  suffice to assure that with probability at least 0.99, any consistent learner using  $H$  will output a hypothesis with true error at most 0.05?

$$|H| = 3^4$$

$$m \geq \frac{1}{0.05} \left( \ln(|H|) + \ln\left(\frac{1}{0.01}\right) \right)$$

### Example: Depth 2 Decision Trees $m \geq \frac{1}{\epsilon}(\ln |H| + \ln(1/\delta))$

Consider classification problem  $f: X \rightarrow Y$ :

- instances:  $X = \langle X_1 \dots X_N \rangle$  where each  $X_i$  is boolean
- learned hypotheses are decision trees of depth 2, using only two variables

How many training examples  $m$  suffice to assure that with probability at least 0.99, *any* learner that outputs a consistent depth 2 decision tree will have true error at most 0.05?

### Example: Depth 2 Decision Trees $m \geq \frac{1}{\epsilon}(\ln |H| + \ln(1/\delta))$

Consider classification problem  $f: X \rightarrow Y$ :

- instances:  $X = \langle X_1 \dots X_N \rangle$  where each  $X_i$  is boolean
- learned hypotheses are decision trees of depth 2, using only two variables

$$\binom{N}{2} \text{ Trees} = \frac{N!}{(N-2)!2!} = \frac{N(N-1)}{2} \quad 2^4 \text{ ways to label the nodes}$$

$$|H| = 8N(N-1)$$

How many training examples  $m$  suffice to assure that with probability at least 0.99, *any* learner that outputs a consistent depth 2 decision tree will have true error at most 0.05?

$$m \geq \frac{1}{0.05} \left( \ln(8N(N-1)) + \ln\left(\frac{1}{0.01}\right) \right)$$

## PAC learning

Consider a class  $C$  of possible target concepts defined over a set of instances  $X$  of length  $n$ , and a learner  $L$  using hypothesis space  $H$ .

*Definition:*  $C$  is **PAC-learnable** by  $L$  using  $H$  if for all  $c \in C$ , distributions  $\mathcal{D}$  over  $X$ ,  $\epsilon$  such that  $0 < \epsilon < 1/2$ , and  $\delta$  such that  $0 < \delta < 1/2$ , learner  $L$  will with probability at least  $(1 - \delta)$  output a hypothesis  $h \in H$  such that  $\text{error}_{\mathcal{D}}(h) \leq \epsilon$ , in time that is polynomial in  $1/\epsilon$ ,  $1/\delta$ ,  $n$  and  $\text{size}(c)$ .

## PAC learning

Consider a class  $C$  of possible target concepts defined over a set of instances  $X$  of length  $n$ , and a learner  $L$  using hypothesis space  $H$ .

*Definition:*  $C$  is **PAC-learnable** by  $L$  using  $H$  if for all  $c \in C$ , distributions  $\mathcal{D}$  over  $X$ ,  $\epsilon$  such that  $0 < \epsilon < 1/2$ , and  $\delta$  such that  $0 < \delta < 1/2$ , learner  $L$  will with probability at least  $(1 - \delta)$  output a hypothesis  $h \in H$  such that  $\text{error}_{\mathcal{D}}(h) \leq \epsilon$ , in time that is polynomial in  $1/\epsilon$ ,  $1/\delta$ ,  $n$  and  $\text{size}(c)$ .

Sufficient condition:  
Holds if learner  $L$  requires only a polynomial number of training examples, and processing per example is polynomial

## Agnostic learning

So far, assumed  $c \in H$

Agnostic learning setting: don't assume  $c \in H$

- What do we want then?
  - The hypothesis  $h$  that makes fewest errors on training data
- What is sample complexity in this case?

$$m \geq \frac{1}{2\epsilon^2}(\ln |H| + \ln(1/\delta))$$

Here  $\epsilon$  is the difference between the training error and true error of the output hypothesis (the one with lowest training error)

## Additive Hoeffding Bounds – Agnostic Learning

- Given  $m$  independent flips of a coin with true  $\Pr(\text{heads}) = \theta$  we can bound the error  $\epsilon$  in the maximum likelihood estimate  $\hat{\theta}$

$$\Pr[\theta > \hat{\theta} + \epsilon] \leq e^{-2m\epsilon^2}$$

- Relevance to agnostic learning: for any single hypothesis  $h$

$$\Pr[\text{error}_{\text{true}}(h) > \text{error}_{\text{train}}(h) + \epsilon] \leq e^{-2m\epsilon^2}$$

- But we must consider all hypotheses in  $H$

$$\Pr[(\exists h \in H) \text{error}_{\text{true}}(h) > \text{error}_{\text{train}}(h) + \epsilon] \leq |H|e^{-2m\epsilon^2}$$

- So, with probability at least  $(1-\delta)$  every  $h$  satisfies

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2m}}$$

## General Hoeffding Bounds

- When estimating parameter  $\theta$  inside  $[a,b]$  from  $m$  examples

$$P(|\hat{\theta} - E[\hat{\theta}]| > \epsilon) \leq 2e^{\frac{-2m\epsilon^2}{(b-a)^2}}$$

- When estimating a probability  $\theta$  is inside  $[0,1]$ , so

$$P(|\hat{\theta} - E[\hat{\theta}]| > \epsilon) \leq 2e^{-2m\epsilon^2}$$

- And if we're interested in only one-sided error, then

$$P((E[\hat{\theta}] - \hat{\theta}) > \epsilon) \leq e^{-2m\epsilon^2}$$

$$m \geq \frac{1}{2\epsilon^2}(\ln |H| + \ln(1/\delta))$$

Here  $\epsilon$  is the difference between the training error and true error of the output hypothesis (this holds for all  $h$  in  $H$ )

But, the output  $h$  with lowest training error might not give us the  $h^*$  with lowest true error. How far can true error of  $h$  be from  $h^*$ ?



$$m \geq \frac{1}{2\epsilon^2}(\ln |H| + \ln(1/\delta))$$

Here  $\epsilon$  is the difference between the training error and true error of the output hypothesis (this holds for all  $h$  in  $H$ )

But, the output  $h$  with lowest training error might not give us the  $h^*$  with lowest true error. How far can true error of  $h$  be from  $h^*$  ?

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{true}}(h^*) + 2\epsilon$$

best training error  
hypothesis

best true error  
hypothesis

$$m \geq \frac{1}{\epsilon}(\ln |H| + \ln(1/\delta))$$

Question: If  $H = \{h \mid h: X \rightarrow Y\}$  is infinite, what measure of complexity should we use in place of  $|H|$  ?

$$m \geq \frac{1}{\epsilon}(\ln |H| + \ln(1/\delta))$$

Question: If  $H = \{h \mid h: X \rightarrow Y\}$  is infinite, what measure of complexity should we use in place of  $|H|$  ?

Answer: The largest subset of  $X$  for which  $H$  can guarantee zero training error (regardless of how it is labeled)

$$m \geq \frac{1}{\epsilon}(\ln |H| + \ln(1/\delta))$$

Question: If  $H = \{h \mid h: X \rightarrow Y\}$  is infinite, what measure of complexity should we use in place of  $|H|$  ?

Answer: The largest subset of  $X$  for which  $H$  can guarantee zero training error (regardless of the target function  $c$ )

**VC dimension of  $H$  is the size of this subset**

Question: If  $H = \{h \mid h: X \rightarrow Y\}$  is infinite, what measure of complexity should we use in place of  $|H|$  ?

Answer: The largest subset of  $X$  for which  $H$  can guarantee zero training error (regardless of the target function  $c$ )

Informal intuition:

- decision tree example: how many labels do we need to see to learn  $h$ ?

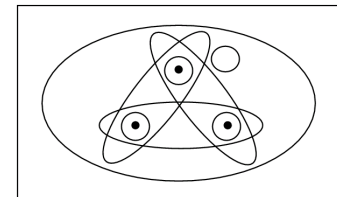
## Shattering a set of instances

*Definition:* a **dichotomy** of a set  $S$  is a partition of  $S$  into two disjoint subsets.

a labeling of each member of  $S$  as positive or negative

*Definition:* a set of instances  $S$  is **shattered** by hypothesis space  $H$  if and only if for every dichotomy of  $S$  there exists some hypothesis in  $H$  consistent with this dichotomy.

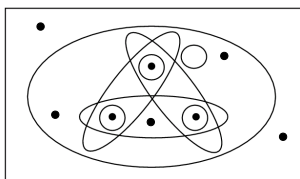
Instance space  $X$



## The Vapnik-Chervonenkis Dimension

*Definition:* The **Vapnik-Chervonenkis dimension**,  $VC(H)$ , of hypothesis space  $H$  defined over instance space  $X$  is the size of the largest finite subset of  $X$  shattered by  $H$ . If arbitrarily large finite sets of  $X$  can be shattered by  $H$ , then  $VC(H) \equiv \infty$ .

Instance space  $X$



$VC(H)=3$

## Sample Complexity based on VC dimension

How many randomly drawn examples suffice to  $\epsilon$ -exhaust  $VS_{H,D}$  with probability at least  $(1-\delta)$ ?

ie., to guarantee that any hypothesis that perfectly fits the training data is probably  $(1-\delta)$  approximately  $(\epsilon)$  correct

$$m \geq \frac{1}{\epsilon} (4 \log_2(2/\delta) + 8VC(H) \log_2(13/\epsilon))$$

Compare to our earlier results based on  $|H|$ :

$$m \geq \frac{1}{\epsilon} (\ln(1/\delta) + \ln |H|)$$

## VC dimension: examples

Consider  $X = \mathbb{R}$ , want to learn  $c: X \rightarrow \{0,1\}$

What is VC dimension of



- Open intervals:

H1: if  $x > a$  then  $y = 1$  else  $y = 0$

H2: if  $x > a$  then  $y = 1$  else  $y = 0$   
or, if  $x > b$  then  $y = 0$  else  $y = 1$

- Closed intervals:

H3: if  $a < x < b$  then  $y = 1$  else  $y = 0$

H4: if  $a < x < b$  then  $y = 1$  else  $y = 0$   
or, if  $a < x < b$  then  $y = 0$  else  $y = 1$

## VC dimension: examples

Consider  $X = \mathbb{R}$ , want to learn  $c: X \rightarrow \{0,1\}$

What is VC dimension of



- Open intervals:

H1: if  $x > a$  then  $y = 1$  else  $y = 0$       VC(H1)=1

H2: if  $x > a$  then  $y = 1$  else  $y = 0$       VC(H2)=2  
or, if  $x > b$  then  $y = 0$  else  $y = 1$

- Closed intervals:

H3: if  $a < x < b$  then  $y = 1$  else  $y = 0$       VC(H3)=2

H4: if  $a < x < b$  then  $y = 1$  else  $y = 0$       VC(H4)=3  
or, if  $a < x < b$  then  $y = 0$  else  $y = 1$

## VC dimension: examples

What is VC dimension of lines in a plane?

- $H_2 = \{ ((w_0 + w_1x_1 + w_2x_2) > 0 \rightarrow y=1) \}$



## VC dimension: examples

What is VC dimension of

- $H_2 = \{ ((w_0 + w_1x_1 + w_2x_2) > 0 \rightarrow y=1) \}$   
VC(H<sub>2</sub>)=3

- For  $H_n$  = linear separating hyperplanes in  $n$  dimensions,  
VC(H<sub>n</sub>)=n+1



For any finite hypothesis space  $H$ , can you give an upper bound on  $VC(H)$  in terms of  $|H|$  ?  
(hint: yes)

## More VC Dimension Examples to Think About

- Logistic regression over  $n$  continuous features
  - Over  $n$  boolean features?
- Decision trees defined over  $n$  boolean features  
 $F: \langle X_1, \dots, X_n \rangle \rightarrow Y$
- Decision trees of depth 2 defined over  $n$  features
- Naïve Bayes defined over  $n$  boolean features
- How about 1-nearest neighbor?

## Tightness of Bounds on Sample Complexity

How many examples  $m$  suffice to assure that any hypothesis that fits the training data perfectly is probably  $(1-\delta)$  approximately  $(\epsilon)$  correct?

$$m \geq \frac{1}{\epsilon} (4 \log_2(2/\delta) + 8VC(H) \log_2(13/\epsilon))$$

How tight is this bound?

## Shatter coefficient $H[m]$

for  $S \subseteq X$ , where  $S = \{x_1 \dots x_m\}$ , define  $H(S)$  as the set of distinct labelings of  $S$  induced by  $H$

$$H(S) \equiv \{ \langle h(x_1) \dots, h(x_m) \rangle \mid h \in H \}$$

and define  $H[m]$  as the maximum number of ways to label  $m$  instances of  $X$

$$H[m] \equiv \max_{S \subseteq X, |S|=m} |H(S)|$$

If  $H$  can shatter a subset of size  $m$ , then  $H[m] = 2^m$

Note  $VCdim(H) \equiv$  largest  $m$  for which  $H[m] = 2^m$

## Shatter coefficient $H[m]$

**Sauer's Lemma:** Let  $VCdim(H) = d$ . Then

1. for all  $m$ ,  $H[m] \leq \Phi_d(m)$ , where  $\Phi_d(m) \equiv \sum_{i=0}^d \binom{m}{i}$
2. for  $m > d$ ,

$$\Phi_d(m) \leq (1 + m)^d$$

$$\Phi_d(m) \leq \left(\frac{em}{d}\right)^d$$

## Sample Complexity - Summary

How many randomly drawn examples suffice to  $\epsilon$ -exhaust  $VS_{H,D}$  with probability at least  $(1-\delta)$ ?

ie., to guarantee that any hypothesis that perfectly fits the training data is probably  $(1-\delta)$  approximately  $(\epsilon)$  correct

$$m \geq \frac{1}{\epsilon} (\ln(1/\delta) + \ln |H|)$$

$|H|$

$$m \geq \frac{1}{\epsilon} (4 \log_2(2/\delta) + 8VC(H) \log_2(13/\epsilon))$$

$VC(H)$

$$m > \frac{2}{\epsilon} (\log_2(1/\delta) + \log_2(3 H[2m]))$$

$H[m]$

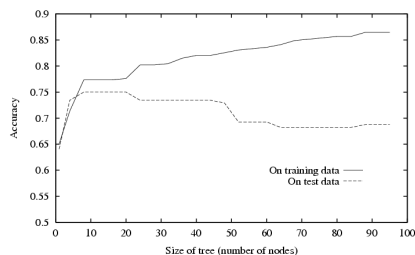
\* also Rademacher complexity

## Agnostic Learning: VC Bounds

[Schölkopf and Smola, 2002]

With probability at least  $(1-\delta)$  every  $h \in H$  satisfies

$$error_{true}(h) < error_{train}(h) + \sqrt{\frac{VC(H)(\ln \frac{2m}{VC(H)} + 1) + \ln \frac{4}{\delta}}{m}}$$

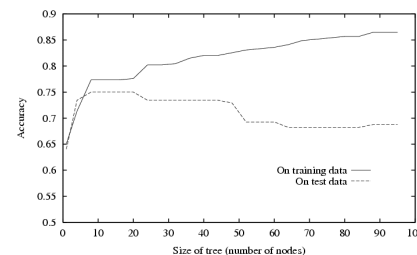


## Agnostic Learning: VC Bounds

[Schölkopf and Smola, 2002]

With probability at least  $(1-\delta)$  every  $h \in H$  satisfies

$$error_{true}(h) < error_{train}(h) + \sqrt{\frac{VC(H)(\ln \frac{2m}{VC(H)} + 1) + \ln \frac{4}{\delta}}{m}}$$



## Sample Complexity - Summary

How many randomly drawn examples suffice to  $\epsilon$ -exhaust  $VS_{H,D}$  with probability at least  $(1-\delta)$ ?

ie., to guarantee that any hypothesis that perfectly fits the training data is probably  $(1-\delta)$  approximately  $(\epsilon)$  correct

$$m \geq \frac{1}{\epsilon} (\ln(1/\delta) + \ln |H|) \quad |H|$$

$$m \geq \frac{1}{\epsilon} (4 \log_2(2/\delta) + 8VC(H) \log_2(13/\epsilon)) \quad VC(H)$$

$$m > \frac{2}{\epsilon} (\log_2(1/\delta) + \log_2(3 H[2m])) \quad H[m]$$

\* also Rademacher complexity

With probability  $\geq (1 - \delta)$ ,  $(error_{true} - error_{train}) \leq \epsilon$

(1) for all  $h \in H$  such that  $error_{train} = 0$ ,

$$\epsilon = \frac{\ln |H| + \ln(1/\delta)}{m} \quad \text{finite } H$$

(2) for all  $h \in H$

$$\epsilon = \sqrt{\frac{\ln |H| + \ln(1/\delta)}{2m}} \quad \text{Agnostic} \quad \text{finite } H$$

(3) for all  $h \in H$

$$\epsilon = 8 \sqrt{\frac{VC(H)(\ln \frac{m}{VC(H)} + 1) + \ln(8/\delta)}{2m}} \quad \text{Agnostic} \quad \text{infinite } H$$

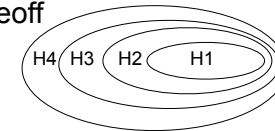
## We stopped here in lecture 14

## Structural Risk Minimization

[Vapnik]

Which hypothesis space should we choose?

- Bias / variance tradeoff



SRM: choose  $H$  to minimize bound on expected true error!

$$error_{true}(h) < error_{train}(h) + \sqrt{\frac{VC(H)(\ln \frac{2m}{VC(H)} + 1) + \ln \frac{4}{\delta}}{m}}$$

\* unfortunately a somewhat loose bound...

## Rademacher Complexity

Key idea: complexity of  $H$  is its ability to fit noise labels.

Advantages:

- applies to real-valued functions (e.g., regression)
- is sensitive to  $P(X)$ , and particular training set
- gives tighter bounds than VC dimension
- widely used in modern learning theory

## Rademacher Complexity Setting

Learn  $f : X \rightarrow Y$ , where  $Y \in \{-1, +1\}$

Note:

if  $h(x) = y$ , then  $yh(x) = 1$

if  $h(x) \neq y$ , then  $yh(x) = -1$

so error of  $h$  on sample  $S = \{\langle x_1, y_1 \rangle, \dots, \langle x_m, y_m \rangle\}$  is :

$$error_S(h) = \frac{1}{m} \sum_{i=1}^m \delta(h(x_i) \neq y_i) = \frac{1}{m} \sum_{i=1}^m \frac{1 - y_i h(x_i)}{2}$$

and the hypothesis  $h$  with the lowest  $error_S(h)$  is

$$\arg \max_{h \in H} \frac{1}{m} \sum_{i=1}^m y_i h(x_i)$$

## Rademacher complexity

Given data sample  $S = \{\langle x_1, y_1 \rangle, \dots, \langle x_m, y_m \rangle\}$

define corresponding set of random labels  $\{\sigma_1, \dots, \sigma_m\}$

where  $\sigma_i \in \{-1, 1\}$ ,  $P(\sigma_i = -1) = 0.5 = P(\sigma_i = 1)$ .

Note the hypothesis  $h$  that best fits these random labels is

$$\arg \max_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i)$$

Define *empirical Rademacher complexity*  $\hat{R}_S(H)$  with respect to  $S$ :

$$\hat{R}_S(H) \equiv E_{\sigma} \left[ \max_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right]$$

## Rademacher complexity

Given data sample  $S = \{\langle x_1, y_1 \rangle, \dots, \langle x_m, y_m \rangle\}$

define corresponding set of random labels  $\{\sigma_1, \dots, \sigma_m\}$

where  $\sigma_i \in \{-1, 1\}$ ,  $P(\sigma_i = -1) = 0.5 = P(\sigma_i = 1)$ .

Note the hypothesis  $h$  that best fits these random labels is

$$\arg \max_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i)$$

Define *empirical Rademacher complexity*  $\hat{R}_S(H)$  with respect to  $S$ :

$$\hat{R}_S(H) \equiv E_{\sigma} \left[ \max_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right]$$

then in the agnostic PAC learning setting, with probability  $(1 - \delta)$ :

$$error_{true}(h) \leq error_{train}(h) + \hat{R}_{train}(H) + 3\sqrt{\frac{\log(2/\delta)}{m}}$$

## Rademacher complexity

$$\hat{R}_S(H) \equiv E_\sigma \left[ \max_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right]$$

What is  $\hat{R}_S(H)$  when:

$H = \{h_1\}$  has only one hypothesis?

$H$  can shatter the training set  $S$ ?

## Rademacher complexity

$$\hat{R}_S(H) \equiv E_\sigma \left[ \max_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right]$$

What is  $\hat{R}_S(H)$  when:

$H = \{h_1\}$  has only one hypothesis?  $\hat{R}_S(H) = 0$

$H$  can shatter the training set  $S$ ?  $\hat{R}_S(H) = 1$

## Empirical Rademacher Complexity

$$\hat{R}_S(H) \equiv E_\sigma \left[ \max_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right]$$

Rademacher complexity:

- applies to real-valued functions (e.g., regression)
- is sensitive to  $P(X)$ , and particular training set
- can give tighter bounds than VC dimension

Also define full *Rademacher complexity*

$$R_m(H) \equiv E_{S \text{ of size } m} [\hat{R}_S(H)]$$

With probability  $\geq (1 - \delta)$ ,  $(error_{true} - error_{train}) \leq \epsilon$

(1) for all  $h \in H$  such that  $error_{train} = 0$ ,

$$\epsilon = \frac{\ln |H| + \ln(1/\delta)}{m}$$

finite H

(2) for all  $h \in H$

$$\epsilon = \sqrt{\frac{\ln |H| + \ln(1/\delta)}{2m}}$$

Agnostic

finite H

(3) for all  $h \in H$

$$\epsilon = 8 \sqrt{\frac{VC(H) \left( \ln \frac{m}{VC(H)} + 1 \right) + \ln(8/\delta)}{2m}}$$

Agnostic

infinite H

(4) for all  $h \in H$

$$\epsilon = \hat{R}_{train}(H) + 3 \sqrt{\frac{\log(2/\delta)}{m}}$$

Agnostic

infinite H